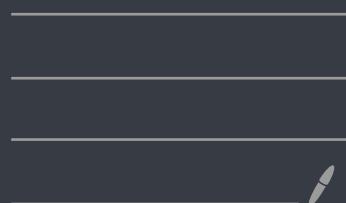


Machine learning

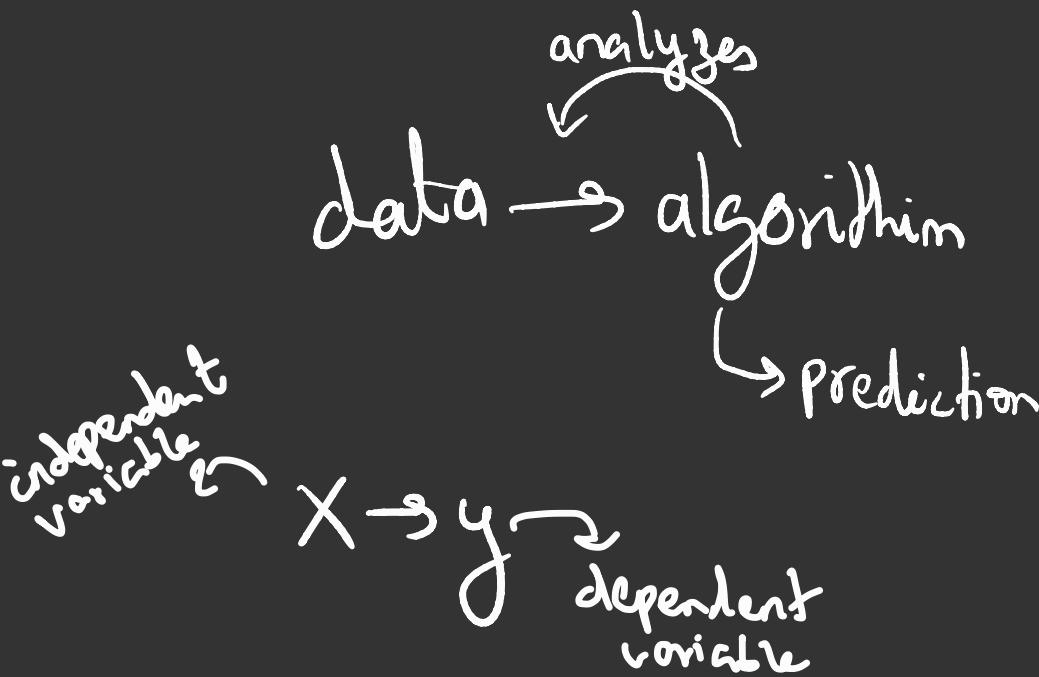
freeCodeCamp



What is Machine learning?

Ans- Computer programs that uses algorithms to analyze data and make intelligent predictions based on the data without being explicitly programmed

Data is given to the algorithm which analyzes it and gives out predictions



X is the input feature

y is the output variable

$$f(n) \rightarrow y$$

$$f(n_1, n_2, n_3) \rightarrow y$$

eg - Email spam detection system

$$T \rightarrow 0/1 \quad (0, \text{is not spam}) \\ 1, \text{is spam}$$

E → Experience

P → Performance

* A computer program is said to learn from experience E with respect to some task T and some performance measure P,

improves with experience E

How it works?

Ans) ★ Study the problem

★ Train the algorithm -

★ Evaluate it

★ If it is good

* Else, Analyze the error

e.g. Spam detection

→ Study the problem

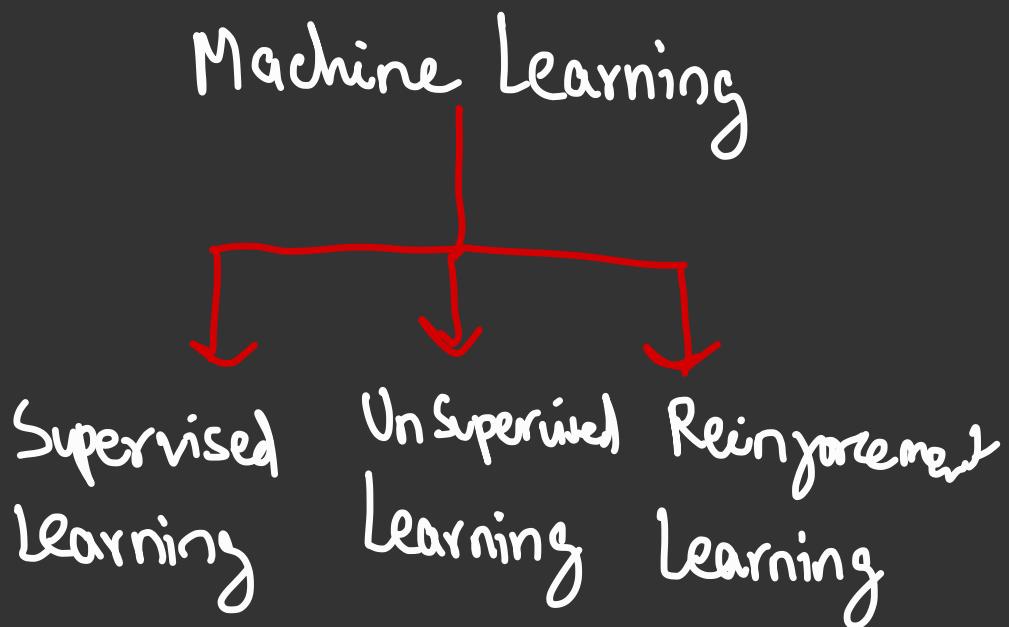
→ $f(n)$ → algorithm

→ evaluate

ML is
a very
iterative
process

- Launch the system
- Error Analysis

Types of Machine Learning Systems

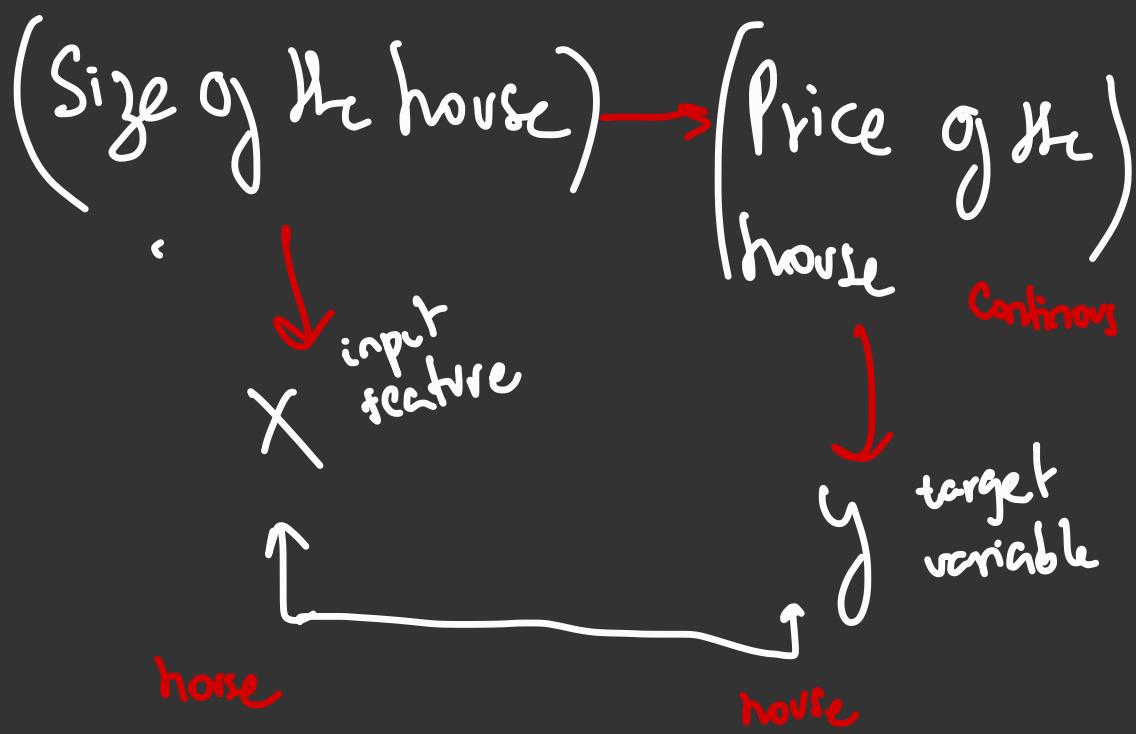


Supervised Learning:

In this the data we feed is properly organised with proper labels, such as → CSV, excel files etc.

In which the data are labeled and we know what our output should be having the relationship b/w input values "X" and output values "y"

eg- House price prediction



in Supervised Learning we're given

$$\{(x_i, y_i)\}_{i=1}^m \rightarrow \text{where } m \text{ is the data.}$$

Unsupervised Learning :-

In this the data we get is not organised nor properly labelled; And we don't know what our output should like and there is not any kind of relationship between input variables and output variables, we have to recognize patterns based on the data.

Eg - working on data like

images, pdfs, audio etc.

$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_m, y_m)\}$$

Types of Supervised Learning problems :-

Supervised Learning



Regression

Classification

Regression - In this type of Supervised learning we're supposed to predict a specific number from a continuous data.

Eg: Predicting House prices

Classification: In this type of supervised learning we're supposed to predict categories of data from categorical data

Eg → Predict is it a apple/orange, cat/dog, yes/no problems

Dividing our data :-

We divide our data into 2/3

sets \Rightarrow Train - test - validation

Train - test

Train Set - The sample of data used to fit the model and train it

Validation Set - The sample of data used to provide an unbiased evaluation of a model fit on the training set

While tuning model hyper parameters.

The evaluation becomes more biased as skill on the validation dataset is incorporated into the model configuration.

Test Dataset: The sample of data used to provide an unbiased evaluation of a final model fit on the training data set.

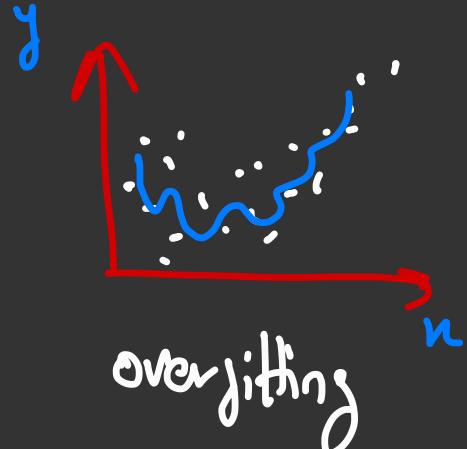
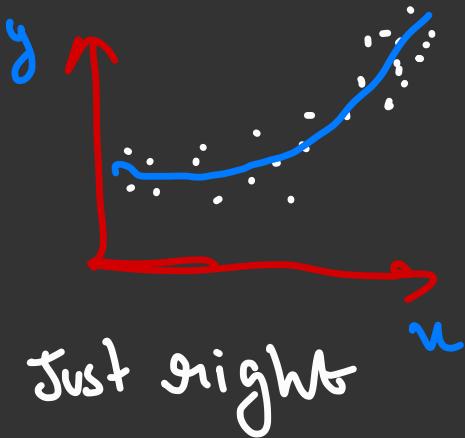
* Sometimes people use test-train-validation set to be extra cautious

- ★ Other times people use only train-test sets using test as both validation and test.
- ★ There is no optimal split percentage but the common split percentage include:
 - Train: 80%, Test: 20%.
 - Train: 70%, Test: 30%.
 - Train: 80%, Test: 10%, validation: 10%.
 - Train: 70%, Test: 15%, validation: 15%.

Overshitting the training data:

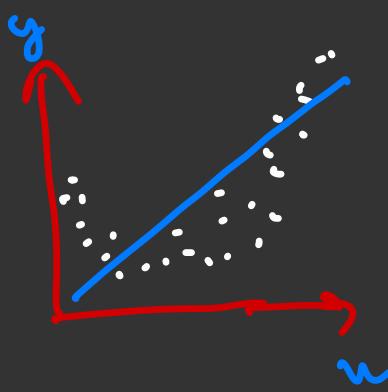
It means the model performs well on the training data, but it does not generalize well.

Solution to this problem is to gather more data, to reduce noise in the training data, to use regularisation.

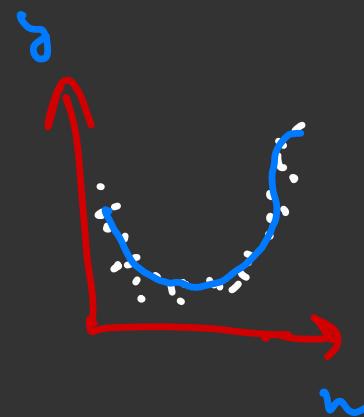


Underfitting the training Data:

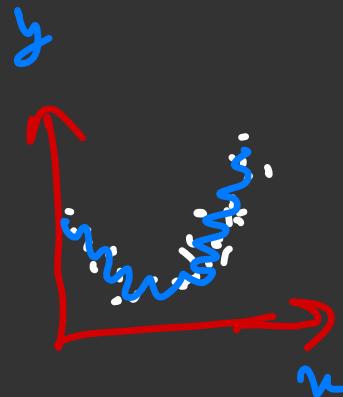
It means the model performs bad on the training data, so it's obvious it will also perform badly on testing data.



underfitted



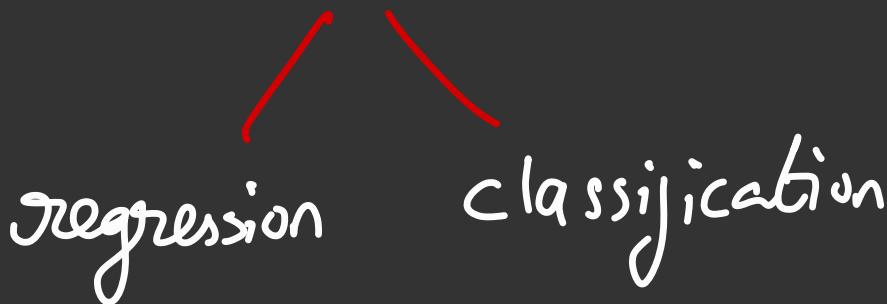
good fit
Robust



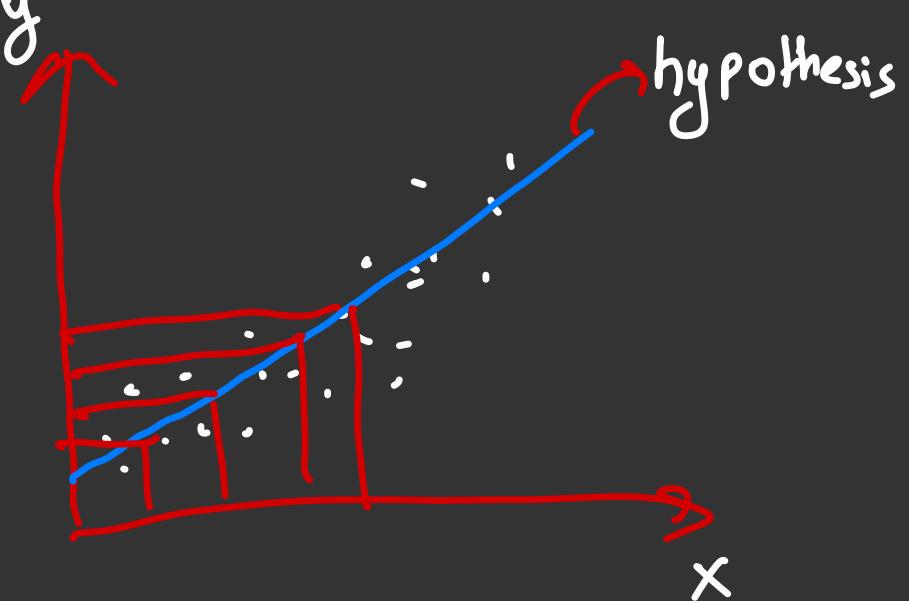
Overfitted

Linear Regression :-

Supervised



linear regression



hypothesis is a straight line as near to all the data points

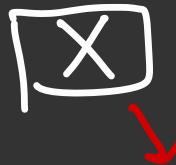
hypothesis function :

$$\boxed{y} = \theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 \dots + \theta_n \cdot x_n$$

* x_0 is always 1



the y intercept



Every X is a different
input feature

$$\therefore f(n) = \theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 \dots + \theta_n \cdot x_n$$

$$= \boxed{y}$$

Eg- We have to predict the price
of house from size of the house and
no. of jars in it.

let $\theta_0 = 2$

$\theta_1 = 3$

$\theta_2 = 4$

$\therefore \Rightarrow 2 + 3 \cdot \text{size of house} + 4 \cdot \text{no. of fns}$

$$= \begin{bmatrix} y \\ \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

We have the input features from the data, we only need to figure out θ for each input feature.

$$f(x) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n$$

Vector form:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

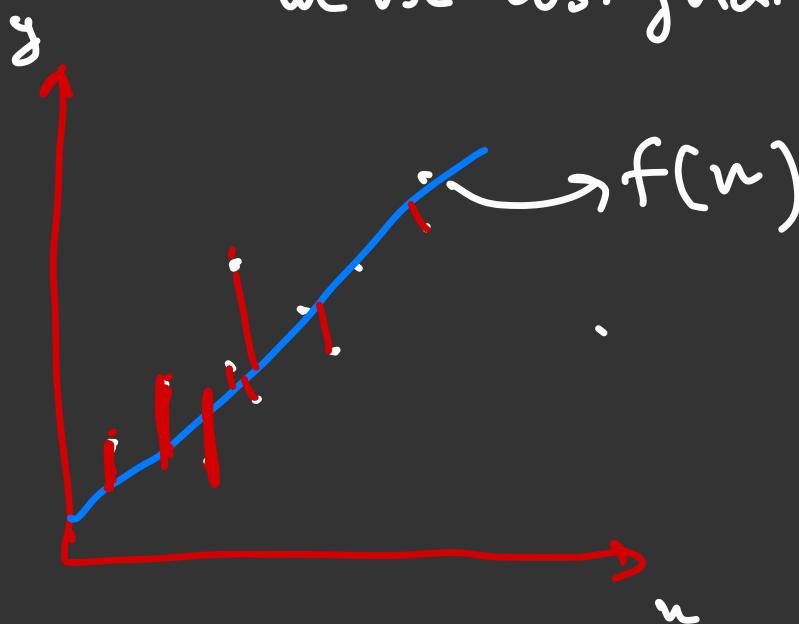
$$f(n) = \theta^T \cdot X$$

(np.dot(theta, n))

or

$$f(n) = \sum_{i=0} \theta_i^T \cdot X_i$$

To evaluate your model :
we use cost function



We find the distance b/w predicted and actual values.

i.e \Rightarrow Predicted value - actual value

$$J(\theta) = \frac{1}{m} + \sum_{i=1}^m (y - \hat{y})^2$$

MSE
(mean squared)
Error

No. of data points

$$f(x) = \sum_{i=0}^n \theta_i x_i$$

actual values

Predicted value

* is square root $J(\theta)$ i.e $\sqrt{J(\theta)}$

it is called RMSE.

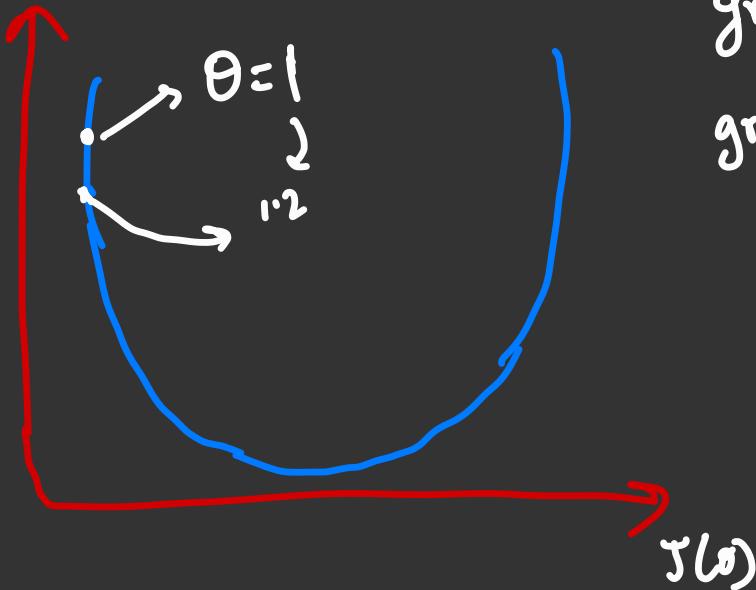
(root mean squared error).

higher RMSE or MSE, the badder
is the model

How to find the most optimial θ ?

We use gradient descent algorithm
which is also called optimization
algorithm)

graphing
gradient
descent
algorithm



it keeps tweaking the θ
values until the cost function is 0

$$J(\theta) \approx 0$$

to do that we find out the
 $\frac{\partial}{\partial \theta} J(\theta) \Rightarrow$ partial derivative
of $J(\theta)$

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{2}{m} + \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

$\theta := \theta - \alpha \frac{\partial}{\partial \theta} (J(\theta))$

Algorithm for gradient descent

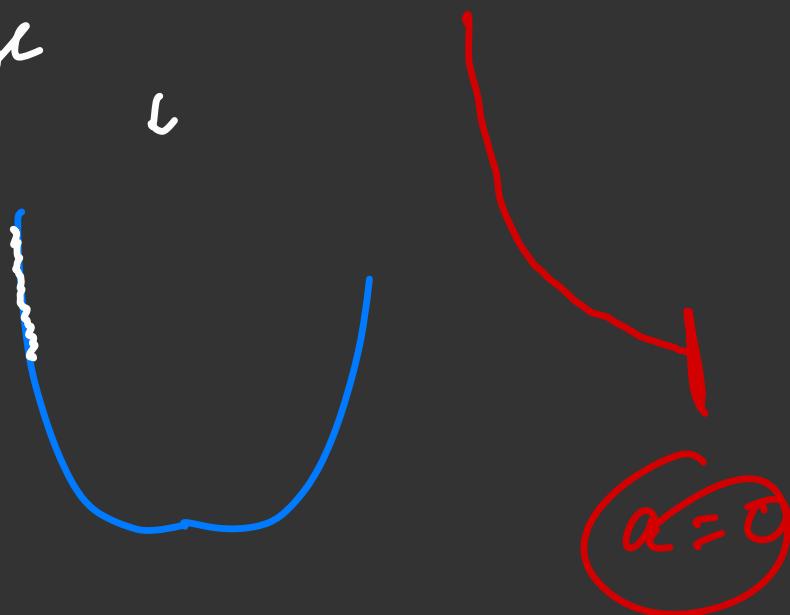
new θ old θ Learning rate old θ

* if α is too large it will never converge

L



* if θ is too small : it will never diverge



* Optimal α :



$$\nabla J = 0$$

- 0.1 for large dataset
- 0.01 for medium to smaller
- 0.001 for smaller
- 0.0001 for very small dataset

You can tune it using

Random CV

Vector form -

$$\frac{\partial}{\partial} J(\theta) = \begin{bmatrix} \frac{\partial}{\partial} J(\theta_0) \\ \frac{\partial}{\partial} J(\theta_1) \\ \vdots \\ \frac{\partial}{\partial} J(\theta_n) \end{bmatrix}$$

Normal
Equation

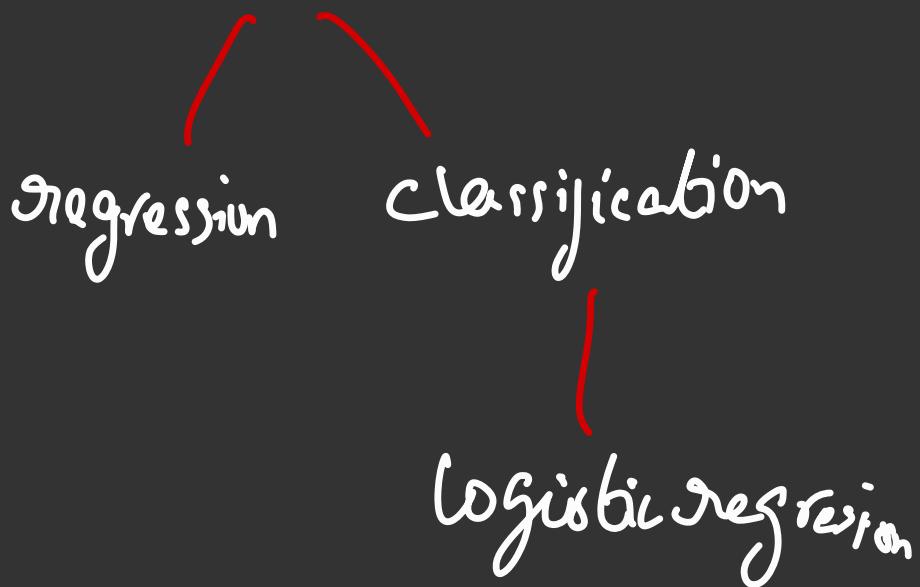
$$\theta = (X^T \cdot X)^{-1} X^T y$$

Polynomial Regression :-

in this we transform our data
to be fitting over linearly

Logistic Regression

Supervised



$$h(u) = \theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n$$



$$\hat{y} = \sigma(h(u))$$

Sigmoid function



restrict the value
b/w 0 and 1

$$\therefore \sigma(h(u)) = \frac{1}{1+e^{-h(u)}}$$

Cost function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h(x^{(i)}) + (1-y^{(i)}) \log(1-h(x^{(i)}))$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h(x^{(i)}) + (1-y^{(i)}) \log(1-h(x^{(i)}))$$

$$h(u) = \sigma(h(u))$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - y^{(i)}) \cdot x_j$$

$$\theta_j = \theta_i - \alpha \frac{\partial}{\partial} J(\theta)$$



Linear regression from scratch:

Approximation:

$$\hat{y} = w_n + b$$

$w \Rightarrow$ slope

$b \Rightarrow$ bias/intercept

$\hat{y} \Rightarrow$ predicted
values

$x \Rightarrow$ data

Cost function:

$$MSE = J(\omega, b) =$$

$$\frac{1}{N} \sum_{i=1}^N (y_i - (\omega_i + b))^2$$

MSE \Rightarrow Mean Squared Error

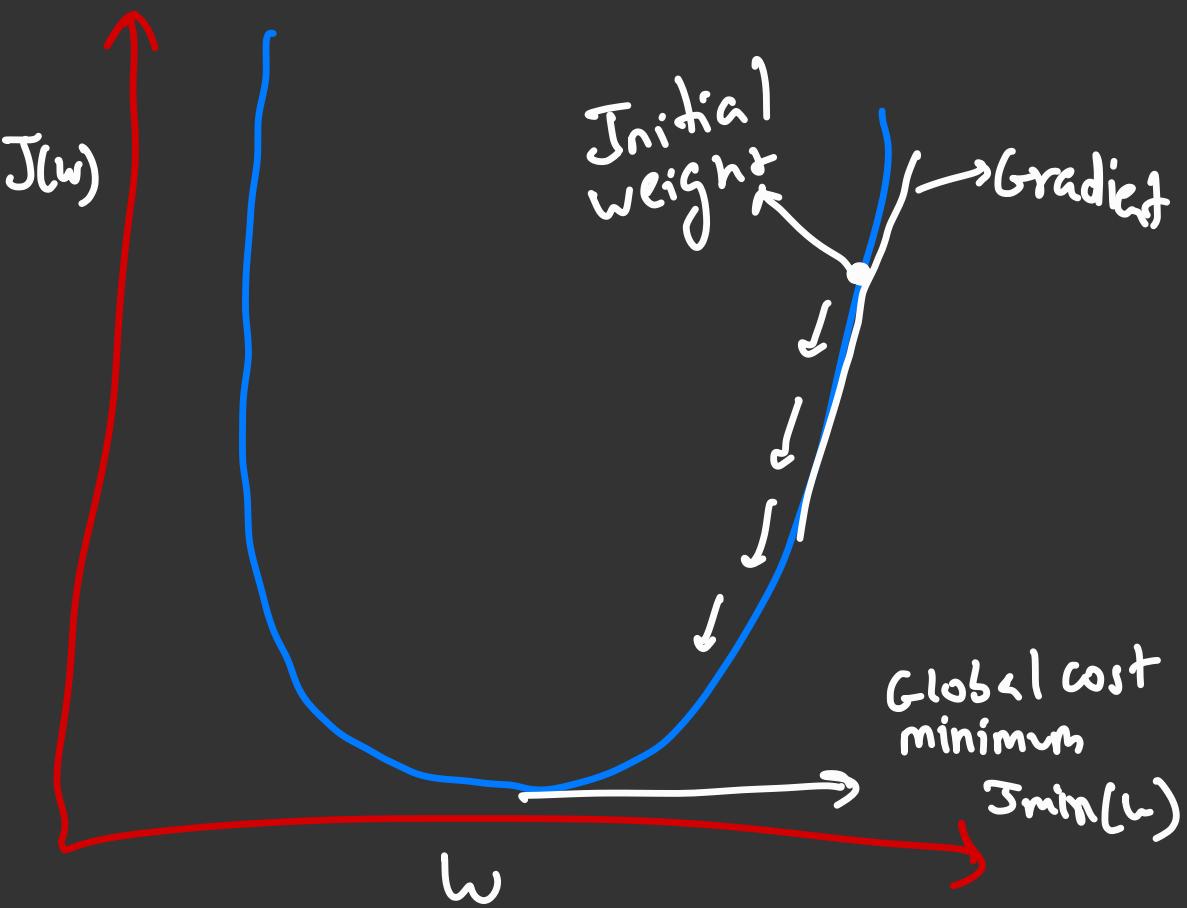
N \Rightarrow Total length of data

y_i \Rightarrow actual values

We use gradient descent to find the w and b values which gives the lowest MSE

$$\nabla J(w, b) = \begin{bmatrix} \frac{\partial f}{\partial w} \\ \frac{\partial f}{\partial b} \end{bmatrix}$$

$$= \left[\frac{1}{n} \sum -2y_i (y_i - (w_i + b)) \right] \\ \left[\frac{1}{n} \sum -2 (y_i - (w_i + b)) \right]$$



(same graph b/w $J(b)$ and b)

Updated values become

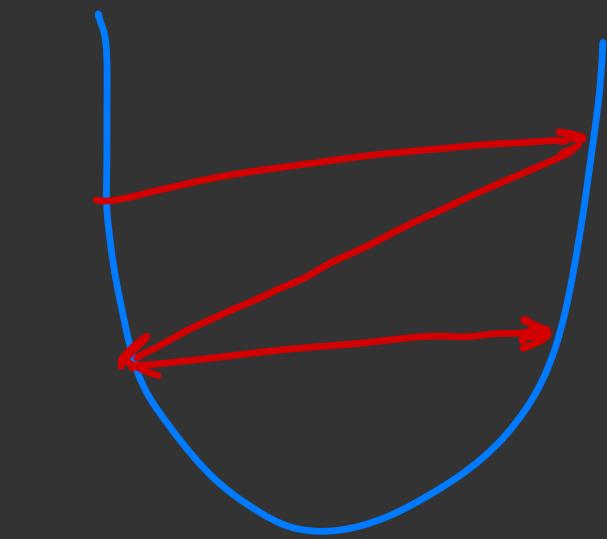
$$w = w - \alpha \cdot dw$$

$$b = b - \alpha \cdot db$$

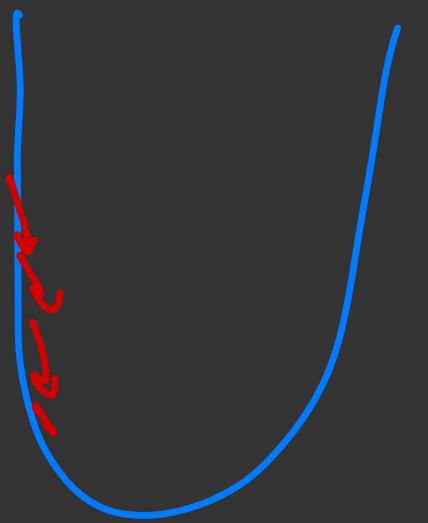
$\alpha \Rightarrow$ Learning rate

$$\frac{dJ}{dw} \approx dw = \frac{1}{N} \sum_{i=1}^N -2n_i (y_i - (w n_i + b))$$

$$\frac{dJ}{db} \approx db = \frac{1}{N} \sum_{i=1}^N -2(y_i - (w n_i + b))$$



big learning
rate



small learning
rate